

Lehrstuhl für
Rechnerarchitektur & Parallele Systeme
Prof. Dr. Martin Schulz
Dominic Prinz
Jakob Schäffeler

Lehrstuhl für
Design Automation
Prof. Dr.-Ing. Robert Wille
Stefan Engels

Einführung in die Rechnerarchitektur

Wintersemester 2024/2025

Übungsblatt 5: Caching

18.11.2024 – 22.11.2024

1 Speicherzugriffszeit

Gegeben sei ein Programm mit $3 \cdot 10^9$ Instruktionen, das Sie auf einem in-order execution Prozessor mit einer Frequenz von 1.0 GHz ausführen. Dieses Programm läuft auf diesem Prozessor mit durchschnittlich 2 Zyklen pro Instruktion (CPI), exklusive Speicherzugriffe. Jede Instruktion generiere dabei durchschnittlich 1.33 Speicherzugriffe.

- a) Wie groß ist die durchschnittliche Speicherzugriffszeit (Average Memory Access Time), wenn das Programm eine Ausführungszeit von 66 Sekunden hat?

$$CPUTIME = IC \cdot \left(\frac{CPI}{f} + \frac{\text{Memory Accesses}}{\text{Instr}} \cdot AMAT \right)$$

$$66 \text{ s} = 3 \cdot 10^9 \text{ instr} \cdot \left(\frac{2 \frac{\text{cycles}}{\text{instr}}}{1 \cdot 10^9 \frac{\text{cycles}}{\text{s}}} + 1.33 \frac{\text{MA}}{\text{instr}} \cdot AMAT \right)$$

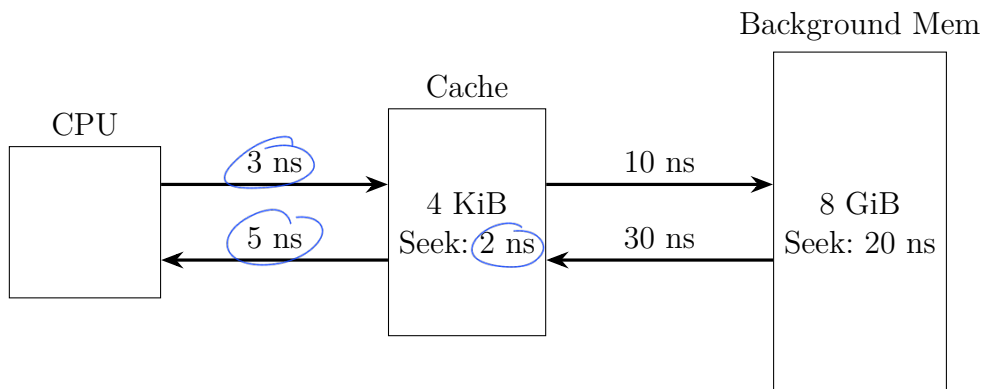
$$66 \text{ s} = \frac{6 \cdot 10^9 \text{ cycles}}{1 \cdot 10^9 \frac{\text{cycles}}{\text{s}}} + 4 \cdot 10^9 \text{ MA} \cdot AMAT$$

$$66 \text{ s} = 6 \text{ s} + 4 \cdot 10^9 \text{ MA} \cdot AMAT$$

$$60 \text{ s} = 4 \cdot 10^9 \text{ MA} \cdot AMAT$$

$$AMAT = \frac{60 \text{ s}}{4 \cdot 10^9 \text{ MA}} = 15 \cdot 10^{-9} \frac{\text{s}}{\text{MA}} = 15 \frac{\text{ns}}{\text{MA}}$$

Der Speicher eines Computers sei wie folgt aufgebaut:



- b) Berechnen Sie Hit-Latency, Miss Penalty, Miss-Latency und Miss-Rate.

Hit-Latency: $3 \text{ ns} + 2 \text{ ns} + 5 \text{ ns} = 10 \text{ ns}$

Miss-Penalty: 60 ns

Miss-Latency: $3 \text{ ns} + 2 \text{ ns} + 10 \text{ ns} + 20 \text{ ns} + 30 \text{ ns} + 5 \text{ ns} = 70 \text{ ns}$

- c) Ein*e Programmierer*in optimiert das Programm um die Hit-Rate im Cache zu maximieren. Dadurch senkt sich die durchschnittliche Speicherzugriffszeit um $\frac{1}{3} = 33.33\%$. Dennoch erhöht sich die Ausführungszeit. Woran könnte das liegen?

- Mehr Instanzen
- Mehr Speicherzugriffe
- Mehr Zyklen/Instanzen

- d) Angenommen neben der um $\frac{1}{3}$ reduzierten durchschnittlichen Speicherzugriffszeit verändert sich durch die Optimierung nur die Anzahl an Instruktionen. Wie sehr muss sich die Anzahl an Instruktionen relativ verändern, damit sich die Laufzeit nicht verbessert?

Ansatz: Stellen Sie Formel für CPU Time nach IC um und setzen Sie einmal die reguläre Average Memory Access Time und dann die verkürzte Average Memory Access Time ein. Um die relative Abweichung zu berechnen, teilen Sie dann die beiden Instruction Counts durcheinander.

1) c)

$$AMAT = \text{Hit Rate} \cdot \text{Hit Latency} + \text{Miss Rate} \cdot \text{Miss Latency}$$

$$AMAT = (1 - \text{Missrate}) \cdot \text{Hit Latency} + \text{Miss Rate} \cdot \text{Miss Latency}$$

$$19 \text{ ns} = (1 - \text{Missrate}) \cdot 10 \text{ ns} + \text{Miss Rate} \cdot 70 \text{ ns}$$

$$19 \text{ ns} = 10 \text{ ns} - MR \cdot 10 \text{ ns} + \text{Missrate} \cdot 70 \text{ ns}$$

$$9 \text{ ns} = -MR \cdot 10 \text{ ns} + \text{Missrate} \cdot 70 \text{ ns}$$

$$9 \text{ ns} = MR \cdot 60 \text{ ns}$$

$$MR = \frac{9 \text{ ns}}{60 \text{ ns}} = \frac{1}{12} \approx 8,3 \%$$

1) d)

$$K = \frac{CPUTIME}{\left(\frac{CPI}{f} + \frac{MA}{f_{ns}} \cdot AMAT \right)}$$

$$\frac{K_2}{K_1} = \frac{\frac{CPUTIME}{\left(\frac{CPI}{f} + \frac{MA}{f_{ns}} \cdot AMAT_2 \right)}}{\frac{CPUTIME}{\left(\frac{CPI}{f} + \frac{MA}{f_{ns}} \cdot AMAT_1 \right)}}$$

$$\frac{K_2}{K_1} = \frac{\left(\frac{CPI}{f} + \frac{MA}{f_{ns}} \cdot AMAT_1 \right)}{\left(\frac{CPI}{f} + \frac{MA}{f_{ns}} \cdot AMAT_2 \right)}$$

$$= \frac{\left(\frac{2 \frac{\text{cycles}}{ns}}{10^9 \frac{\text{cycles}}{s}} + 1,33 \frac{MA}{ns} \cdot 19 \cdot 10^{-9} \frac{s}{ns} \right)}{\left(\frac{2 \frac{\text{cycles}}{ns}}{10^9 \frac{\text{cycles}}{s}} + 1,33 \frac{MA}{ns} \cdot 10 \cdot 10^{-9} \frac{s}{ns} \right)}$$

$$= \frac{\left(2 \cdot 10^{-9} \frac{s}{ns} + 1,33 \frac{MA}{ns} + 19 \cdot 10^{-9} \frac{s}{ns} \right)}{\left(2 \cdot 10^{-9} \frac{s}{ns} + 1,33 \frac{MA}{ns} + 10 \cdot 10^{-9} \frac{s}{ns} \right)}$$

$$= \frac{2 + 1,33 \cdot 19}{2 + 1,33 \cdot 10} = \frac{21,95 \frac{s}{ns}}{19,3 \frac{s}{ns}} = 1,134$$

e) Was sind Gründe, weshalb sich die Laufzeiten eines Programms heutzutage nicht mehr so einfach wie oben berechnen lassen?

- Parallelisierung: mehrere Prozesse (z.B. DS), mehrere Rechenkern
- Verschieden Cache-Elemente
- Out-of-Order-Prozessierung, Sortierung, etc.

2 Vergleich der Cachearten

Gegeben sei eine 16-Bit Architektur mit byte-adressierbarem Speicher. Zudem ist ein 256 Byte großer Cache mit 2 Byte langen Cachezeilen gegeben. Vergleichen Sie in den folgenden Teilaufgaben die verschiedenen Cache-Typen, indem Sie die Fragen jeweils für einen direct mapped Cache, einen 4-fach assoziativen Cache und einen voll-assoziativen Cache beantworten.

a) In wie viele unterschiedliche Cachezeilen kann ein Byte aus dem Speicher potentiell gecached werden?

$$256 \text{ Byte} / 2 \text{ Byte} = 128 \text{ Cachezeilen}$$

Voll-assoziativ: 128

4-fach: 4

Direct mapped: 1

$$\# \text{ Cachezeilen} = \frac{128}{4} = 32$$

b) Es wird auf das Byte an der Adresse 0xABCD zugegriffen. Wie viele unterschiedliche Cachezeilen müssen überprüft werden, um festzustellen, ob das Datum an der Adresse bereits gecached ist? Geben Sie zudem an, welcher Vergleich stattfindet, um festzustellen, ob der Wert bereits gecached ist.

0xABCD = 0b 1010 1011 1100 1101

Voll-assoziativ: Offset-Bits: 1

Index-Bits: 0

Tag-Bits: 16 - 1 = 15 1010 1011 1100 1101

4-fach: Offset-Bits: 1

Index-Bits: 5

Tag-Bits: 16 - 1 - 5 = 10 1010 1011 11

Direct mapped: Offset-Bits: 1

Index-Bits: 7

Tag-Bits: 16 - 1 - 7 = 8 1010 1011

$\log_2 128$

3 Cachestruktur und Cachemisses

Für die folgenden Fragen wird ein System betrachtet, das nur einen Rechenkern mit einem Cache besitzt. Dieser Cache sei 2-fach assoziativ, habe 4 Cachezeilenmengen (Cache Sets) und eine Cachezeilenlänge von 32 Byte. Die Architektur sei eine reine 32-Bit Architektur und alle Zugriffe 32 Bit breit. Der Speicher sei wie üblich Byte-adressierbar.

a) Wie viele Cachezeilen hat der Cache?

$$4 \text{ Cache Sets} * 2 \frac{\text{Cachezeilen}}{\text{Set}} = 8 \text{ Cachezeilen}$$

b) Wie viele Bits benötigt man für Offset, Index und Tag?

Offset: $\log_2 32 = 5 \text{ Bits}$

Index: $\log_2 4 = 2 \text{ Bits}$

Tag: $32 - 5 - 2 = 25 \text{ Bit}$

c) Gegeben sei folgender Ausschnitt eines C-Programms:

```
int data[1000];
int sum = 0;
```

// ABSCHNITT 1

```
for (int i = 0; i < 1000; i += 64) {
    sum = sum + data[i];
}
```

// ABSCHNITT 2

```
for (int i = 0; i < 1000; i += 32) {
    sum = sum + data[i];
}
```



Zugriffe

$$\left\lceil \frac{1000}{64} \right\rceil = 16$$

Zugriffe

$$\left\lceil \frac{1000}{32} \right\rceil = 32$$

Nehmen Sie an, dass nur das Array im Speicher liegt und alle anderen Werte in Registern gehalten werden. Der Cache sei zu Beginn leer. Gehen Sie davon aus, dass ein int 32-Bit groß ist und stets vollständig in einer Cachezeile liegt.

i) Wie viele „Capacity Misses“ sind in den beiden Abschnitten zu erwarten?

Startadresse : 0

Adresse₆₄ : Startadresse +

$$4 \cdot 64 = 256$$

1. Abschnitt: 0 CM

Adresse₃₂ : Startadresse +

$$4 \cdot 32 = 128$$

2. Abschnitt: 16

Iteration	Tag	Index	Offset
1	00000	00	00000
2	0001	00	00000
3	00010	00	00000

ii) Wie viele „Cold Misses“ sind zu erwarten?

Ansatz: Betrachten Sie die Adressen der einzelnen Array-Elemente.

1. Abschnitt: 16 Cold Misses

2. Abschnitt: 16 Cold Misses

~~0, 32, 64, 96, 128~~

- iii) Knobelaufgabe: Wie ändern sich diese Zahlen, wenn die Variablen `sum` und `i` nicht in Registern sondern direkt nach dem Array `data` befinden?

4 Ersetzungsstrategien

Gegeben sei ein 4-fach assoziativer Cache mit 16 Cachezeilen. Pro Cachezeile werden 4 Byte abgespeichert. In den folgenden Tabellen ist jeweils das Cache Set 0 angegeben. Cachen sie die folgenden Werte jeweils einmal mit LFU und LRU als Ersetzungsstrategie und füllen sie in den Tabellen aus. Als Vereinfachung wird von jeder Adresse nur das Tag angegeben und Index und Offset werden in dieser Aufgabe ignoriert. Falls die Ersetzungsstrategie zu keinem eindeutigen Ergebnis führt ersetzen Sie immer die niederwertigere Cacheline mit dem neuen Tag.

- a) Tragen Sie in die Tabelle die für die Blöcke gespeicherten Tags mit **LRU** als Ersetzungsstrategie ein.

*Beispiel: In Zeitschritt 1 wird die Adresse mit dem Tag 12 gecached. Da das Datum an der Adresse noch nicht im Cache vorhanden ist handelt es sich um ein Miss. Der Wert wird in Cachezeile 1 und der letzte Access in **Last Use** von Cachezeile 1 gespeichert.*

Time	Tag	Line 0		Line 1		Line 2		Line 3		Hit/Miss
		Tag	LU	Tag	LU	Tag	LU	Tag	LU	
1	12	12	1	—	—	—	—	—	—	Miss
2	20	12	1	20	2	—	—	—	—	Miss
3	12	12	3	20	2	—	—	—	—	Hit
4	12	12	4	20	2	—	—	—	—	Hit
5	28	12	4	20	2	28	5	—	—	Miss
6	32	12	4	20	2	28	5	32	6	Miss
7	38	12	4	38	7	28	5	32	6	Miss
8	25	25	8	38	7	28	5	32	6	Miss
9	12	25	8	38	7	12	9	32	6	Miss
10	32	25	8	38	7	12	9	32	10	Hit

- b) Tragen Sie in die Tabelle die für die Blöcke gespeicherten Tags mit **LFU** als Ersetzungsstrategie ein.

Beispiel: In Zeitschritt 1 wird die Adresse mit dem Tag 12 gecached. Da das Datum an der Adresse noch nicht im Cache vorhanden ist handelt es sich um ein Miss. Der Wert wird in Cachezeile 1 gespeichert und der Counter auf 1 gesetzt.

Time	Tag	Line 0		Line 1		Line 2		Line 3		Hit/Miss
		Tag	Count	Tag	Count	Tag	Count	Tag	Count	
1	12	12	1	—	—	—	—	—	—	Miss
2	20	12	1	20	1					Miss
3	12	12	2	20	1					Hit
4	12	12	3	20	1					Hit
5	28	12	3	20	1	28	1			Miss
6	32	12	3	20	1	28	1	32	1	Miss
7	38	12	3	38	1	28	1	32	1	Miss
8	25	12	3	25	1	28	1	32	1	Hit
9	12	12	4	25	1	28	1	32	1	Hit
10	32	12	4	25	1	28	1	32	2	Hit

c) Vergleichen Sie die Ersetzungsstrategien. Was sind Vor- und Nachteile der jeweiligen Strategien?

- LFU: Gut repetition Zugriffsmuster / Schlecht bei unregelmäßige Zugriffsmuster
- LRU: Schlecht bei zufälligen / verteilten Mustern

5 Cache-Simulation (Hausaufgabe 5)

Bearbeitung und Abgabe der Hausaufgabe 5 auf <https://artemis.in.tum.de/courses/401> bis Sonntag, den 24.11.2024, 23:59 Uhr.